



## Finite-Sample Analysis of LSTD

Alessandro Lazaric, Mohammad Ghavamzadeh, Remi Munos

### ► To cite this version:

Alessandro Lazaric, Mohammad Ghavamzadeh, Remi Munos. Finite-Sample Analysis of LSTD. ICML - 27th International Conference on Machine Learning, Jun 2010, Haifa, Israel. pp.615-622. inria-00482189

**HAL Id: inria-00482189**

**<https://inria.hal.science/inria-00482189>**

Submitted on 9 May 2010

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

# Finite-Sample Analysis of LSTD

---

Alessandro Lazaric

Mohammad Ghavamzadeh

Rémi Munos

INRIA Lille - Nord Europe, Team SequeL, FRANCE

ALESSANDRO.LAZARIC@INRIA.FR

MOHAMMAD.GHAVAMZADEH@INRIA.FR

REMI.MUNOS@INRIA.FR

## Abstract

In this paper we consider the problem of policy evaluation in reinforcement learning, i.e., learning the value function of a fixed policy, using the least-squares temporal-difference (LSTD) learning algorithm. We report a finite-sample analysis of LSTD. We first derive a bound on the performance of the LSTD solution evaluated at the states generated by the Markov chain and used by the algorithm to learn an estimate of the value function. This result is general in the sense that no assumption is made on the existence of a stationary distribution for the Markov chain. We then derive generalization bounds in the case when the Markov chain possesses a stationary distribution and is  $\beta$ -mixing.

## 1. Introduction

Least-squares temporal-difference (LSTD) learning (Bradtke & Barto, 1996; Boyan, 1999) is a widely used algorithm for prediction in general, and in the context of reinforcement learning (RL), for learning the value function  $V^\pi$  of a given policy  $\pi$ . LSTD has been successfully applied to a number of problems especially after the development of the least-squares policy iteration (LSPI) algorithm (Lagoudakis & Parr, 2003), which extends LSTD to control problems. More precisely, LSTD computes the fixed point of the operator  $\Pi\mathcal{T}$ , where  $\mathcal{T}$  is the Bellman operator and  $\Pi$  is the projection operator in a linear function space. Although LSTD and LSPI have been widely used in the RL community, a finite-sample analysis of LSTD (i.e., performance bounds in terms of the number of samples) is still lacking.

Most of the theoretical work analyzing LSTD have

been focused on the model-based case, where explicit models of the reward function and the dynamics are available. In particular, Tsitsiklis & Van Roy (1997) showed that the distance between the LSTD solution and the value function  $V^\pi$  is bounded by the distance between  $V^\pi$  and its closest approximation in the linear space, multiplied by a constant which increases as the discount factor approaches 1. In this bound, it is assumed that the Markov chain possesses a stationary distribution  $\rho^\pi$  and the distances are measured according to  $\rho^\pi$ . Bertsekas (2001) reported a similar analysis for the empirical version of LSTD. His analysis reveals a critical dependency on the inverse of the smallest eigenvalue of the LSTD's  $A$  matrix (note that the LSTD solution is obtained by solving the system of linear equations  $Ax = b$ ). Nonetheless, Bertsekas (2001) does not provide a finite-sample analysis of the algorithm. On the other hand, Antos et al. (2008) analyzed the modified Bellman residual (MBR) minimization algorithm for a finite number of samples, bounded function spaces, and a  $\mu$ -norm that might be different from the norm induced by  $\rho^\pi$ . Although MBR minimization was shown to reduce to LSTD in case of linear spaces, it is not straightforward how the finite-sample bounds derived by Antos et al. (2008) can be extended to unbounded linear spaces considered by LSTD.

In this paper, we report a finite-sample analysis of LSTD. To the best of our knowledge, this is the first complete finite-sample analysis of this widely used algorithm. Our analysis is for a specific implementation of LSTD that we call *pathwise LSTD*. Pathwise LSTD has two specific characteristics: **1)** it takes a single trajectory generated by the Markov chain induced by policy  $\pi$  as input, and **2)** it uses the pathwise Bellman operator (will be precisely defined later), which is defined to be a contraction w.r.t. the empirical norm. We first derive a bound on the performance of the pathwise LSTD solution for a setting that we call *Markov design*. In this setting, the performance is evaluated at the points used by the algorithm to learn an estimate of  $V^\pi$ . This bound is general in the sense that no as-

sumption is made on the existence of a stationary distribution for the Markov chain. Then, in the case the Markov chain admits a stationary distribution  $\rho^\pi$  and is  $\beta$ -mixing, we derive generalization bounds w.r.t. the norm induced by  $\rho^\pi$ .

Besides providing a full finite-sample analysis of LSTD, the major insights gained by the analysis in the paper can be summarized as follows. The first result is about the existence of the LSTD solution and its performance. In Theorem 1 we show that with a slight modification of the empirical Bellman operator  $\widehat{\mathcal{T}}$  (leading to the definition of pathwise LSTD), the operator  $\widehat{\Pi}\widehat{\mathcal{T}}$  (where  $\widehat{\Pi}$  is an empirical projection operator) has always a fixed point  $\hat{v}$  even when the sample-based Gram matrix is not invertible and the Markov chain does not admit a stationary distribution. In this very general setting, it is still possible to derive a bound for the performance of  $\hat{v}$  evaluated on the states of the trajectory, and an analysis of the bound reveals a critical dependency on the smallest strictly positive eigenvalue  $\nu_n$  of the sample-based Gram matrix. Then, in the case in which the Markov chain has a stationary distribution  $\rho^\pi$ , it is possible to relate the value of  $\nu_n$  to the smallest eigenvalue of the Gram matrix defined according to  $\rho^\pi$ . Furthermore, it is possible to generalize the previous performance bound over the entire state space under the measure  $\rho^\pi$ , when the samples are drawn from a stationary  $\beta$ -mixing process (Theorem 2). It is important to note that the asymptotic bound obtained by taking the number of samples,  $n$ , to infinity is equal (up to constants) to the bound in Tsitsiklis & Van Roy (1997) for model-based LSTD. Furthermore, a comparison with the bounds in Antos et al. (2008) shows that we successfully leverage on the specific setting of LSTD: **1**) the space of functions is linear, and **2**) the distribution used to evaluate the performance is the stationary distribution of the Markov chain induced by the policy. In particular, we obtain a better bound both in terms of estimation error, a rate of order  $O(1/n)$  instead of  $O(1/\sqrt{n})$  for the squared error, and in terms of approximation error, the minimal distance between the value function  $V^\pi$  and the space  $\mathcal{F}$  instead of the inherent Bellman errors of  $\mathcal{F}$ . Finally, the extension in Theorem 3 to the case in which the samples belong to a trajectory generated by a fast mixing Markov chain shows that it is possible to achieve the same performance as in the case of stationary  $\beta$ -mixing processes.

## 2. Preliminaries

For a measurable space with domain  $\mathcal{X}$ , we let  $\mathcal{S}(\mathcal{X})$  and  $\mathcal{B}(\mathcal{X}; L)$  denote the set of probability measures

over  $\mathcal{X}$ , and the space of bounded measurable functions with domain  $\mathcal{X}$  and bound  $0 < L < \infty$ , respectively. For a measure  $\rho \in \mathcal{S}(\mathcal{X})$  and a measurable function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , we define the  $\ell_2(\rho)$ -norm of  $f$ ,  $\|f\|_\rho$ , and for a set of  $n$  states  $X_1, \dots, X_n \in \mathcal{X}$ , we define the empirical norm  $\|f\|_n$  as

$$\|f\|_\rho^2 = \int f(x)^2 \rho(dx) \quad \text{and} \quad \|f\|_n^2 = \frac{1}{n} \sum_{t=1}^n f(X_t)^2.$$

The supremum norm of  $f$ ,  $\|f\|_\infty$ , is defined as  $\|f\|_\infty = \sup_{x \in \mathcal{X}} |f(x)|$ .

We consider the standard RL framework (Sutton & Barto, 1998) in which a learning agent interacts with a stochastic environment by following a policy  $\pi$  and this interaction is modeled as a discrete-time discounted Markov chain (MC). A discounted MC is a tuple  $\mathcal{M}^\pi = \langle \mathcal{X}, R^\pi, P^\pi, \gamma \rangle$ , where the state space  $\mathcal{X}$  is a subset of a Euclidean space, the reward function  $R^\pi : \mathcal{X} \rightarrow \mathbb{R}$  is uniformly bounded by  $R_{\max}$ , the transition kernel  $P^\pi$  is such that for all  $x \in \mathcal{X}$ ,  $P^\pi(\cdot|x)$  is a distribution over  $\mathcal{X}$ , and  $\gamma \in (0, 1)$  is a discount factor. The value function of a policy  $\pi$ ,  $V^\pi$ , is the unique fixed-point of the Bellman operator  $\mathcal{T}^\pi : \mathcal{B}(\mathcal{X}; V_{\max} = \frac{R_{\max}}{1-\gamma}) \rightarrow \mathcal{B}(\mathcal{X}; V_{\max})$  defined by<sup>1</sup>

$$(\mathcal{T}^\pi V)(x) = R^\pi(x) + \gamma \int_{\mathcal{X}} P^\pi(dy|x) V(y).$$

To approximate the value function  $V$ , we use a linear approximation architecture with parameters  $\alpha \in \mathbb{R}^d$  and basis functions  $\varphi_j \in \mathcal{B}(\mathcal{X}; L)$ ,  $j = 1, \dots, d$ . We denote by  $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$ ,  $\phi(\cdot) = (\varphi_1(\cdot), \dots, \varphi_d(\cdot))^\top$  the feature vector, and by  $\mathcal{F}$  the linear function space spanned by the basis functions  $\varphi_j$ . Thus  $\mathcal{F} = \{f_\alpha, \alpha \in \mathbb{R}^d\}$ , where  $f_\alpha(\cdot) = \phi(\cdot)^\top \alpha$ .

Let  $(X_1, \dots, X_n)$  be a sample path (or trajectory) of size  $n$  generated by the Markov chain  $\mathcal{M}$ . Let  $v \in \mathbb{R}^n$  and  $r \in \mathbb{R}^n$  such that  $v_t = V(X_t)$  and  $r_t = R(X_t)$  be the value vector and the reward vector, respectively. Also, let  $\Phi = [\phi(X_1)^\top; \dots; \phi(X_n)^\top]$  be the feature matrix defined at the states, and  $\mathcal{F}_n = \{\Phi\alpha, \alpha \in \mathbb{R}^d\} \subset \mathbb{R}^n$  be the corresponding vector space. We denote by  $\widehat{\Pi} : \mathbb{R}^n \rightarrow \mathcal{F}_n$  the orthogonal projection onto  $\mathcal{F}_n$ , defined as  $\widehat{\Pi}y = \arg \min_{z \in \mathcal{F}_n} \|y - z\|_n$ , where  $\|y\|_n^2 = \frac{1}{n} \sum_{t=1}^n y_t^2$ . Note that the orthogonal projection  $\widehat{\Pi}y$  for any  $y \in \mathbb{R}^n$  exists and is unique.

## 3. Pathwise LSTD

Pathwise LSTD is a version of LSTD which takes as input a single path  $X_1, \dots, X_n$  and returns the fixed-

<sup>1</sup>To simplify the notation, we remove the dependency to the policy  $\pi$  and use  $\mathcal{M}$ ,  $R$ ,  $P$ ,  $V$ , and  $\mathcal{T}$  instead of  $\mathcal{M}^\pi$ ,  $R^\pi$ ,  $P^\pi$ ,  $V^\pi$ , and  $\mathcal{T}^\pi$  throughout the paper.

point of the empirical operator  $\widehat{\Pi}\widehat{\mathcal{T}}$ , where  $\widehat{\mathcal{T}} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is the *pathwise Bellman operator* defined as

$$(\widehat{\mathcal{T}}y)_t = \begin{cases} r_t + \gamma y_{t+1} & 1 \leq t < n, \\ r_t & t = n. \end{cases}$$

Note that by defining the operator  $\widehat{P} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  as  $(\widehat{P}y)_t = y_{t+1}$  for  $1 \leq t < n$  and  $(\widehat{P}y)_n = 0$ , we have  $\widehat{\mathcal{T}}y = r + \gamma\widehat{P}y$ . The motivation for using the pathwise Bellman operator is that it is  $\gamma$ -contraction in  $\ell_2$ -norm, i.e., for any  $y, z \in \mathbb{R}^n$ , we have

$$\|\widehat{\mathcal{T}}y - \widehat{\mathcal{T}}z\|_n^2 = \|\gamma\widehat{P}(y - z)\|_n^2 \leq \gamma^2 \|y - z\|_n^2.$$

Moreover, it can be shown that the orthogonal projection  $\widehat{\Pi}$  is non-expansive: since  $\|\widehat{\Pi}y\|_n^2 = \langle y, \widehat{\Pi}y \rangle_n \leq \|y\|_n \|\widehat{\Pi}y\|_n$ , using the Cauchy-Schwarz inequality we obtain  $\|\widehat{\Pi}y\|_n \leq \|y\|_n$ . Therefore, from Banach fixed point theorem, there exists a unique fixed-point  $\hat{v}$  of the mapping  $\widehat{\Pi}\widehat{\mathcal{T}}$ , i.e.,  $\hat{v} = \widehat{\Pi}\widehat{\mathcal{T}}\hat{v}$ . We call  $\hat{v}$  the *pathwise LSTD solution*. Note that the unicity of  $\hat{v}$  does not imply that there exists a unique parameter  $\hat{\alpha}$  such that  $\hat{v} = \Phi\hat{\alpha}$ .

## 4. Markov Design Bound

**Theorem 1.** *Let  $X_1, \dots, X_n$  be a trajectory of the Markov chain, and  $v, \hat{v} \in \mathbb{R}^n$  be the vectors whose components are the value function and the pathwise LSTD solution at  $\{X_i\}_{i=1}^n$ , respectively. Then with probability  $1 - \delta$ , where the probability is w.r.t. the random trajectory, we have*

$$\|\hat{v} - v\|_n \leq \frac{1}{1 - \gamma} \left[ \|v - \widehat{\Pi}v\|_n + \gamma V_{\max} L \sqrt{\frac{d}{\nu_n}} \left( \sqrt{\frac{8 \log(2d/\delta)}{n}} + \frac{1}{n} \right) \right], \quad (1)$$

where the random variable  $\nu_n$  is the smallest strictly-positive eigenvalue of the sample-based Gram matrix  $\frac{1}{n}\Phi^\top\Phi$ .

**Remark 1** When the eigenvalues of the sample-based Gram matrix  $\frac{1}{n}\Phi^\top\Phi$  are non-zero,  $\Phi^\top\Phi$  is invertible, and thus,  $\widehat{\Pi} = \Phi(\Phi^\top\Phi)^{-1}\Phi^\top$ . In this case, the unicity of  $\hat{v}$  implies the unicity of  $\hat{\alpha}$  since

$$\hat{v} = \Phi\hat{\alpha} \implies \Phi^\top\hat{v} = \Phi^\top\Phi\hat{\alpha} \implies \hat{\alpha} = (\Phi^\top\Phi)^{-1}\Phi^\top\hat{v}. \quad (2)$$

Since  $\hat{v}$  is the unique fixed-point of  $\widehat{\Pi}\widehat{\mathcal{T}}$ , it can be replaced by  $\widehat{\Pi}\widehat{\mathcal{T}}\hat{v}$  in Eq. 2. Using the definitions of  $\widehat{\Pi}$  and  $\widehat{\mathcal{T}}$ , we obtain  $\Phi^\top(I - \gamma\widehat{P})\Phi\hat{\alpha} = \Phi^\top r$ . By defining  $A = \Phi^\top(I - \gamma\widehat{P})\Phi$  and  $b = \Phi^\top r$ ,  $\hat{\alpha}$  can be seen as the unique solution of the  $d \times d$  system of linear equations  $A\alpha = b$ .

**Remark 2** Note that in case there exists a constant  $\nu > 0$ , such that with probability  $1 - \delta'$  all the eigenvalues of the sample-based Gram matrix are lower bounded by  $\nu$ , Eq. 1 (with  $\nu_n$  replaced by  $\nu$ ) holds with probability at least  $1 - (\delta + \delta')$ .

**Remark 3** When the sample-based Gram matrix  $\frac{1}{n}\Phi^\top\Phi$  is not invertible, the unicity of  $\hat{v}$  does not imply the unicity of the solution to the system  $\Phi\alpha = \hat{v}$ . However, since  $\hat{v}$  is the unique fixed point of  $\widehat{\Pi}\widehat{\mathcal{T}}$ , the vector  $\hat{v} - \widehat{\Pi}\hat{v}$  is perpendicular to the space  $\mathcal{F}_n$ , and thus,  $\Phi^\top(\hat{v} - \widehat{\Pi}\hat{v}) = 0$ . By replacing  $\hat{v}$  with  $\Phi\alpha$ , we obtain  $\Phi^\top\Phi\alpha = \Phi^\top(r + \gamma\widehat{P}\Phi\alpha)$  and then  $\Phi^\top(I - \gamma\widehat{P})\Phi\alpha = \Phi^\top r$ . Therefore, we still have the same system of equations  $A\alpha = b$  as in Remark 1, with the exact same  $A$  and  $b$ , but now the system may have many solutions.<sup>2</sup> Among all possible solutions, one may choose the one with minimal norm:  $\hat{\alpha} = A^+b$ , where  $A^+$  is the Moore-Penrose pseudo-inverse of  $A$ .

**Remark 4** Theorem 1 provides a bound without any reference to the stationary distribution of the Markov chain. In fact, the bound of Eq. 1 holds even when the chain does not possess a stationary distribution. For example, consider a Markov chain on the real line where the transitions always move the states to the right, i.e.,  $p(X_{t+1} \in dy | X_t = x) = 0$  for  $y \leq x$ . For simplicity assume that the value function  $V$  is bounded and belongs to  $\mathcal{F}$ . This Markov chain is not recurrent, and thus, does not have a stationary distribution. We also assume that the feature vectors  $\phi(X_1), \dots, \phi(X_n)$  are sufficiently independent, so that the eigenvalues of  $\frac{1}{n}\Phi^\top\Phi$  are greater than  $\nu > 0$ . Then according to Theorem 1, pathwise LSTD is able to estimate the value function at the states at a rate  $O(1/\sqrt{n})$ . This may seem surprising because at each state  $X_t$  the algorithm is only provided with a noisy estimation of the expected value of the next state. However, the estimates are unbiased conditioned on the current state, and we will see in the proof that using a concentration inequality for martingale, pathwise LSTD is able to learn a good estimate of the value function at a state  $X_t$  using noisy pieces of information at other states that may be far away from  $X_t$ . In other words, learning the value function at a given state does not require making an average over many samples close to that state. This implies that LSTD does not require the Markov chain to possess a stationary distribution.

**Remark 5** The most critical part of the bound in Eq. 1 is the inverse dependency on the smallest positive eigenvalue  $\nu_n$ . A similar dependency is shown in the LSTD analysis of Bertsekas (2001). The main

<sup>2</sup>Note that since the fixed point  $\hat{v}$  exists, this system always has at least one solution.

difference is that here we have a more complete finite-sample analysis with an explicit dependency on the number of samples and the other characteristic parameters of the problem. Furthermore, if the Markov chain admits a stationary distribution  $\rho$ , we are able to relate the existence of the LSTD solution to the smallest eigenvalue of the Gram matrix defined according to  $\rho$  (see Section 5.1).

In order to prove Theorem 1, we first introduce the model of regression with *Markov design* and then state and prove a Lemma about this model.

**Definition 1.** *The model of regression with **Markov design** is a regression problem where the data  $(X_t, Y_t)_{1 \leq t \leq n}$  are generated according to the following model:  $X_1, \dots, X_n$  is a sample path generated by a Markov chain,  $Y_t = f(X_t) + \xi_t$ , where  $f$  is the target function, and the noise term  $\xi_t$  is a random variable which is adapted to the filtration generated by  $X_1, \dots, X_{t+1}$  and is such that*

$$|\xi_t| \leq C, \quad \text{and} \quad \mathbb{E}[\xi_t | X_1, \dots, X_t] = 0. \quad (3)$$

**Lemma 1** (Regression bound for the Markov design setting). *We consider the model of regression with Markov design from Definition 1. Let  $\hat{w} \in \mathcal{F}_n$  be the least-squares estimate of the (noisy) values  $Y = \{Y_t\}_1^n$ , i.e.,  $\hat{w} = \hat{\Pi}Y$ , and  $w \in \mathcal{F}_n$  be the least-squares estimate of the (noiseless) values  $Z = \{Z_t\}_1^n = \{f(X_t)\}_1^n$ , i.e.,  $w = \hat{\Pi}Z$ . Then for any  $\delta > 0$ , with probability at least  $1 - \delta$ , where the probability is w.r.t. the random sample path  $X_1, \dots, X_n$ , we have*

$$\|\hat{w} - w\|_n \leq CL \sqrt{\frac{2d \log(2d/\delta)}{n\nu_n}}, \quad (4)$$

where  $\nu_n$  is the smallest strictly-positive eigenvalue of the sample-based Gram matrix  $\frac{1}{n}\Phi^\top \Phi$ .

*Proof of Lemma 1.* We define  $\xi \in \mathbb{R}^n$  to be the vector with components  $\xi_t$ , and  $\hat{\xi} = \hat{w} - w = \hat{\Pi}(Y - Z) = \hat{\Pi}\xi$ . Since the projection is orthogonal we have  $\langle \hat{\xi}, \xi \rangle_n = \|\hat{\xi}\|_n^2$ . Since  $\hat{\xi} \in \mathcal{F}_n$ , there exists at least one  $\alpha \in \mathbb{R}^d$  such that  $\hat{\xi} = \Phi\alpha$ , so by Cauchy-Schwarz inequality we have

$$\begin{aligned} \|\hat{\xi}\|_n^2 &= \langle \hat{\xi}, \xi \rangle_n = \frac{1}{n} \sum_{i=1}^d \alpha_i \sum_{t=1}^n \xi_t \varphi_i(X_t) \\ &\leq \frac{1}{n} \|\alpha\|_2 \left[ \sum_{i=1}^d \left( \sum_{t=1}^n \xi_t \varphi_i(X_t) \right)^2 \right]^{1/2}. \end{aligned} \quad (5)$$

Now among the vectors  $\alpha$  such that  $\hat{\xi} = \Phi\alpha$ , we define  $\hat{\alpha}$  to be the one with minimal  $\ell_2$ -norm, i.e.,  $\hat{\alpha} = \Phi^+ \hat{\xi}$ . Let  $K$  denote the null space of  $\Phi$ , which is also the null space of  $\frac{1}{n}\Phi^\top \Phi$ . Then  $\hat{\alpha}$  can be decomposed as

$\hat{\alpha} = \hat{\alpha}_K + \hat{\alpha}_{K^\perp}$ , where  $\hat{\alpha}_K \in K$  and  $\hat{\alpha}_{K^\perp} \in K^\perp$ , and because the decomposition is orthogonal, we have  $\|\hat{\alpha}\|_2^2 = \|\hat{\alpha}_K\|_2^2 + \|\hat{\alpha}_{K^\perp}\|_2^2$ . Since  $\hat{\alpha}$  is of minimal norm among all the vectors  $\alpha$  such that  $\hat{\xi} = \Phi\alpha$ , its component in  $K$  must be zero, thus  $\hat{\alpha} \in K^\perp$ .

The Gram matrix  $\frac{1}{n}\Phi^\top \Phi$  is positive-semidefinite, thus its eigenvectors corresponding to zero eigenvalues generate  $K$  and the other eigenvectors generate its orthogonal complement  $K^\perp$ . Therefore, from the assumption that the smallest strictly-positive eigenvalue of  $\frac{1}{n}\Phi^\top \Phi$  is  $\nu_n$ , we deduce that since  $\hat{\alpha} \in K^\perp$ ,

$$\|\hat{\xi}\|_n^2 = \frac{1}{n} \hat{\alpha}^\top \Phi^\top \Phi \hat{\alpha} \geq \nu_n \hat{\alpha}^\top \hat{\alpha} = \nu_n \|\hat{\alpha}\|_2^2. \quad (6)$$

By using the result of Eq. 6 in Eq. 5, we obtain

$$\|\hat{\xi}\|_n \leq \frac{1}{n\sqrt{\nu_n}} \left[ \sum_{i=1}^d \left( \sum_{t=1}^n \xi_t \varphi_i(X_t) \right)^2 \right]^{1/2}. \quad (7)$$

Now, from Eq. 3, we have that

$$\mathbb{E}[\xi_t \varphi_i(X_t) | X_1, \dots, X_t] = \varphi_i(X_t) \mathbb{E}[\xi_t | X_1, \dots, X_t] = 0,$$

thus  $\xi_t \varphi_i(X_t)$  is a martingale difference sequence w.r.t. the filtration generated by the Markov chain, and one may apply Azuma's inequality to deduce that with probability  $1 - \delta$ ,

$$\left| \sum_{t=1}^n \xi_t \varphi_i(X_t) \right| \leq CL \sqrt{2n \log(2/\delta)}.$$

By a union bound over all features, we have that with probability  $1 - \delta$ , for all  $i = 1 \dots d$ ,

$$\left| \sum_{t=1}^n \xi_t \varphi_i(X_t) \right| \leq CL \sqrt{2n \log(2d/\delta)}. \quad (8)$$

The results follows by combining Eq. 8 with Eq. 7.  $\square$

**Remark about Lemma 1** In the Markov design model considered in this lemma, states  $\{X_t\}_1^n$  are random variables generated according to the Markov chain and the noise terms  $\xi_t$  may depend on the next state  $X_{t+1}$  (but should be centered conditioned on the past  $X_1, \dots, X_t$ ). This lemma will be used in order to prove Theorem 1, where we replace the target function  $f$  with the value function  $V$ , and the noise term  $\xi_t$  with the temporal difference  $r(X_t) + \gamma V(X_{t+1}) - V(X_t)$ .

Note that this lemma is an extension of the bound for the model of regression with deterministic design in which the states  $\{X_t\}_1^n$  are fixed and the noise terms,  $\xi_t$ 's, are independent. In the setting of deterministic design, usual concentration results provide high probability bounds similar to Eq. 4, but without the dependence on  $\nu_n$ . An open question is whether it is



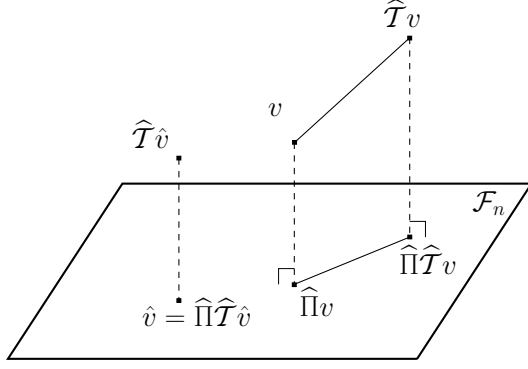


Figure 1. This figure represents the space  $\mathbb{R}^n$ , the linear vector subspace  $\mathcal{F}_n$  and some vectors used in the proof of Theorem 1.

possible to remove  $\nu_n$  in the bound for the Markov design regression setting.

*Proof of Theorem 1. Step 1:* Using the triangle inequality, we have (see Figure 1)

$$\|\hat{v} - v\|_n \leq \|\hat{v} - \hat{\Pi}\hat{\mathcal{T}}v\|_n + \|\hat{\Pi}\hat{\mathcal{T}}v - \hat{\Pi}v\|_n + \|\hat{\Pi}v - v\|_n. \quad (9)$$

From the  $\gamma$ -contraction of  $\widehat{\Pi}\widehat{\mathcal{T}}$  mapping and the fact that  $\hat{v}$  is its unique fixed point, we obtain

$$\|\hat{v} - \hat{\Pi}\hat{\mathcal{T}}v\|_n = \|\hat{\Pi}\hat{\mathcal{T}}\hat{v} - \hat{\Pi}\hat{\mathcal{T}}v\|_n \leq \gamma\|\hat{v} - v\|_n, \quad (10)$$

Thus from Eq. 9 and 10, we have

$$\|\hat{v} - v\|_n \leq \frac{1}{1 - \gamma} \left[ \|\hat{\Pi}v - v\|_n + \|\hat{\Pi}\hat{\mathcal{T}}v - \hat{\Pi}v\|_n \right]. \quad (11)$$

**Step 2:** We now provide a high probability bound on  $\|\widehat{\Pi}\widehat{T}v - \widehat{\Pi}v\|_n$ . This is a consequence of Lemma 1 applied to the vectors  $Y = \widehat{T}v$  and  $Z = v$ . Since  $v$  is the value function at the points  $\{X_t\}_1^n$ , from the definition of the pathwise Bellman operator, we have that for  $1 \leq t \leq n-1$ ,

$$\begin{aligned}\xi_t &= y_t - v_t = r(X_t) + \gamma V(X_{t+1}) - V(X_t) \\ &= \gamma [V(X_{t+1}) - \int P(dy|X_t)V(y)],\end{aligned}$$

and  $\xi_n = y_n - v_n = -\gamma \int P(dy|X_n)V(y)$ . Thus, Eq. 3 holds for  $1 \leq t \leq n-1$ . Here we may choose  $C = 2\gamma V_{\max}$  for a bound on  $\xi_t$ ,  $1 \leq t \leq n-1$ , and  $C = \gamma V_{\max}$  for a bound on  $\xi_n$ . Azuma's inequality may only be applied to the sequence of  $n-1$  terms (the  $n$ -th

term adds a contribution to the bound), thus instead of Eq. 8, we obtain with probability  $1 - \delta$

$$\left| \sum_{t=1}^n \xi_t \varphi_i(X_t) \right| \leq \gamma V_{\max} L (2\sqrt{2n \log(2d/\delta)} + 1),$$

for all  $1 \leq i \leq d$ . Combining with Eq. 7, we deduce that with probability  $1 - \delta$ , we have

$$\|\widehat{\Pi}\widehat{\mathcal{T}}v - \widehat{\Pi}v\|_n \leq \gamma V_{\max} L \sqrt{\frac{d}{\nu_n}} \left( \sqrt{\frac{8 \log(2d/\delta)}{n}} + \frac{1}{n} \right), \quad (12)$$

The claim follows by combining Eq. 12 and 11.  $\square$

**Remark 6** In addition to Eq. 1, one may easily deduce a tighter bound (when  $\gamma$  is close to 1):

$$\begin{aligned} \|\hat{v} - v\|_n &\leq \frac{1}{\sqrt{1 - \gamma^2}} \|v - \hat{\Pi} v\|_n \\ &\quad + \frac{1}{1 - \gamma} \left[ \gamma V_{\max} L \sqrt{\frac{d}{\nu_n}} \left( \sqrt{\frac{8 \log(2d/\delta)}{n}} + \frac{1}{n} \right) \right] \end{aligned}$$

by using Pytagora's Theorem in Step 1, i.e.,  $\|\hat{v}-v\|_n^2 \leq (\|\hat{v}-\hat{\Pi}\hat{\mathcal{T}}v\|_n + \|\hat{\Pi}\hat{\mathcal{T}}v-\hat{\Pi}v\|_n)^2 + \|\hat{\Pi}v-v\|_n^2$  instead of Eq. 9.

## 5. Generalization Bounds

The generality of Theorem 1 comes at the cost that the performance is evaluated only at the states visited by the Markov chain. The reason is that no assumption about the existence of the stationary distribution of the Markov chain is made. However in many problems of interest, the Markov chain has a stationary distribution  $\rho$ , and thus, the performance can be generalized to the whole state space under the measure  $\rho$ . Moreover, if  $\rho$  exists, it is possible to derive a condition for the existence of the pathwise LSTD solution depending on the number of samples and the smallest eigenvalue of the Gram matrix defined according to the stationary distribution  $\rho$ ;  $G \in \mathbb{R}^{d \times d}$ ,  $G_{ij} = \int \phi_i(x) \phi_j(x) \rho(dx)$ . In this section, we assume that the Markov chain  $\mathcal{M}$  is exponentially fast  $\beta$ -mixing with parameters  $\bar{\beta}, b, \kappa$ , i.e., its  $\beta$ -mixing coefficients satisfy  $\beta_i \leq \bar{\beta} \exp(-bi^\kappa)$  (see the appendix for a more detailed definition of  $\beta$ -mixing processes).

Before stating the main results of this section, we introduce some notation. If  $\rho$  is the stationary distribution of the Markov chain, we define the orthogonal projection operator  $\Pi : \mathcal{B}(\mathcal{X}; V_{\max}) \rightarrow \mathcal{F}$  as

$$\Pi V = \arg \min_{f \in \mathcal{F}} \|V - f\|_{\rho}. \quad (13)$$

Furthermore, in the rest of the paper with a little abuse of notation, we replace the empirical norm  $\|v\|_n$  defined on states  $X_1^n$  by  $\|V\|_n$ , where  $V \in \mathcal{B}(\mathcal{X}; V_{\max})$  is such that  $V(X_t) = v_t$ . Finally, we should guarantee that the pathwise LSTD solution  $\hat{V}$  is uniformly bounded on  $\mathcal{X}$ . For this reason, we move from  $\mathcal{F}$  to the truncated space  $\tilde{\mathcal{F}}$ . A function  $\tilde{f} \in \tilde{\mathcal{F}}$  is defined as

$$\tilde{f}(x) = \begin{cases} f(x) & \text{if } |f(x)| \leq V_{\max}, \\ \text{sgn}(f(x))V_{\max} & \text{otherwise.} \end{cases} \quad (14)$$

In the next sections, we present conditions on the existence of the pathwise LSTD solution and derive generalization bounds under different assumptions on the way that the samples  $X_1, \dots, X_n$  are generated.

### 5.1. Existence of Pathwise LSTD Solution

In this section, we assume that all the eigenvalues of  $G$  are strictly positive and derive a condition to guarantee that the sample-based Gram matrix  $\frac{1}{n}\Phi^\top\Phi$  is invertible. In particular, we show that if a large enough number of samples (depending on the smallest eigenvalue of  $G$ ) is available, then the smallest eigenvalue of  $\frac{1}{n}\Phi^\top\Phi$  is strictly positive with high probability.

**Lemma 2.** *Let  $G$  be the Gram matrix defined according to the distribution  $\rho$  and  $\omega > 0$  be its smallest eigenvalue. Let  $X_1, \dots, X_n$  be a path of length  $n$  of a stationary  $\beta$ -mixing process with stationary distribution  $\rho$ . If the number of samples  $n$  satisfies the following condition*

$$\frac{\Lambda(n, \delta)}{n} \max \left\{ \frac{\Lambda(n, \delta)}{b}, 1 \right\}^{1/\kappa} < \frac{\omega}{288L^2}, \quad (15)$$

where  $\Lambda(n, \delta) = \log \frac{e}{\delta} + \log(\max\{6, n\bar{\beta}\})$ , then with probability  $1 - \delta$ , the family of features  $(\varphi_1, \dots, \varphi_d)$  is linearly independent on the states  $X_1, \dots, X_n$  (i.e.,  $\|f_\alpha\|_n = 0$  implies  $\alpha = 0$ ) and the smallest eigenvalue  $\nu_n$  of the sample-based Gram matrix  $\frac{1}{n}\Phi^\top\Phi$  satisfies

$$\sqrt{\nu_n} \geq \frac{\sqrt{\omega}}{2} - \sqrt{72L^2 \frac{\Lambda(n, \delta)}{n} \max \left\{ \frac{\Lambda(n, \delta)}{b}, 1 \right\}^{1/\kappa}} > 0. \quad (16)$$

*Proof.* From the definition of the Gram matrix and the fact that  $\omega$  is its smallest eigenvalue, for any function  $f_\alpha \in \mathcal{F}$ , we have

$$\|f_\alpha\|_\rho^2 = \|\phi^\top \alpha\|_\rho^2 = \alpha^\top G \alpha \geq \omega \alpha^\top \alpha = \omega \|\alpha\|^2.$$

Using the concentration inequality from Corollary 4 in the appendix and the fact that the basis functions  $\varphi_j$  are bounded by  $L$ , thus  $f_\alpha$  is bounded by  $L\|\alpha\|$ , we have  $\|f_\alpha\|_\rho - 2\|f_\alpha\|_n \leq \epsilon$  with probability  $1 - \delta$ , where

$$\epsilon = \|\alpha\| \sqrt{288L^2 \frac{\Lambda(n, \delta)}{n} \max \left\{ \frac{\Lambda(n, \delta)}{b}, 1 \right\}^{1/\kappa}}.$$

Thus we obtain

$$2\|f_\alpha\|_n + \epsilon \geq \sqrt{\omega}\|\alpha\|. \quad (17)$$

Let  $\alpha$  be such that  $\|f_\alpha\|_n = 0$ , then from Eq. 17 and the definition of  $\epsilon$  we deduce that  $\alpha = 0$ . Thus  $\nu_n > 0$  and the inequality in Eq. 16 is obtained by choosing  $\alpha$  to be the eigenvector of  $\frac{1}{n}\Phi^\top\Phi$  correspond to the smallest eigenvalue  $\nu_n$ . For this value of  $\alpha$ , we have  $\|f_\alpha\|_n = \sqrt{\nu_n}\|\alpha\|$ , and the claim follows using Eq. 17.  $\square$

**Remark 1** If  $\Lambda(n, \delta)/b > 1$  and  $n\bar{\beta} \geq 6$ , the condition on the number of samples can be rewritten as

$$\frac{n}{\log \left( \frac{e}{\delta} n\bar{\beta} \right)^{\frac{1+\kappa}{\kappa}}} \geq \frac{288L^2}{\omega b^{1/\kappa}}.$$

As it can be seen, the number of samples needed to have strictly positive eigenvalues in the sample-based Gram matrix has an inverse dependency on the smallest eigenvalue of  $G$ . As a consequence, the more  $G$  is ill-conditioned the more samples we need for the sample-based Gram matrix  $\frac{1}{n}\Phi^\top\Phi$  to be invertible.

### 5.2. Generalization Bounds for Stationary $\beta$ -mixing Processes

In this section, we show how Theorem 1 can be generalized to the entire state space  $\mathcal{X}$  when the Markov chain  $\mathcal{M}$  has a stationary distribution  $\rho$ . In particular, we consider the case in which the samples  $X_1, \dots, X_n$  are obtained by following a single trajectory in the stationary regime of  $\mathcal{M}$ , i.e., when we consider that  $X_1$  is drawn from  $\rho$ .

**Theorem 2.** *Let  $X_1, \dots, X_n$  be a path generated by a stationary  $\beta$ -mixing process with stationary distribution  $\rho$ . Let  $\omega$  be the smallest eigenvalue of the Gram matrix defined according to  $\rho$  and  $n$  satisfy the condition in Eq. 15. Let  $\tilde{V}$  be the truncation (using Eq. 14) of the pathwise LSTD solution, then*

$$\begin{aligned} \|\tilde{V} - V\|_\rho &\leq \frac{2}{1-\gamma} \left[ 2\sqrt{2}\|V - \Pi V\|_\rho + \varepsilon_2 \right. \\ &\quad \left. + \gamma V_{\max} L \sqrt{\frac{d}{\nu}} \left( \sqrt{\frac{8 \log(8d/\delta)}{n}} + \frac{1}{n} \right) \right] + \varepsilon_1 \end{aligned} \quad (18)$$

with probability  $1 - \delta$ , where  $\nu$  is a lower bound on the eigenvalues of the sample-based Gram matrix defined by Eq. 16,

$$\varepsilon_1 = \sqrt{\frac{\Lambda(n, d, \delta/4)}{nC_2} \max \left\{ \frac{\Lambda(n, d, \delta/4)}{b}, 1 \right\}^{1/\kappa}}$$

with  $\Lambda(n, d, \delta/4) = 2(d+1) \log n + \log \frac{4e}{\delta} + \log^+ (\max\{18(C_1 C_2)^{2(d+1)}, \bar{\beta}\})$ ,  $C_1 = 6912eV_{\max}^2$ , and  $C_2 = (1152V_{\max}^2)^{-1}$ , and

$$\varepsilon_2 = \sqrt{288(V_{\max} + L\|\alpha^*\|)^2 \frac{\Lambda(n, \delta/4)}{n} \max\left\{\frac{\Lambda(n, \delta/4)}{b}, 1\right\}}^{1/\kappa}$$

where  $\Lambda(n, \delta/4) = \log \frac{4e}{\delta} + \log(\max\{6, n\bar{\beta}\})$  and  $\alpha^*$  is such that  $f_{\alpha^*} = \Pi V$ .

*Proof.* This result is a consequence of applying generalization bounds to both sides of Eq. 1 (Theorem 1). We first bound the left-hand side.

$$2\|\hat{V} - V\|_n \geq 2\|\tilde{V} - V\|_n \geq \|\tilde{V} - V\|_\rho - \varepsilon_1$$

with probability  $1 - \delta'$ . The first step follows from the definition of the truncation operator, while the second step is a straightforward application of Corollary 3 in the appendix.

We now bound the term  $\|V - \hat{\Pi}V\|_n$  in Eq. 1:

$$\|V - \hat{\Pi}V\|_n \leq \|V - \Pi V\|_n \leq 2\sqrt{2}\|V - \Pi V\|_\rho + \varepsilon_2$$

with probability  $1 - \delta'$ . The first step follows from the definition of the operator  $\hat{\Pi}$ . The second step is an application of the inequality of Corollary 4 in the appendix for the function  $V - \Pi V$ .

From Theorem 1, the two generalization bounds, and the lower bound on  $\nu$ , each one holding with probability  $1 - \delta'$ , the statement of the Theorem (Eq. 18) holds with probability  $1 - \delta$  by setting  $\delta = 4\delta'$ .  $\square$

**Remark 1** Rewriting the bound in terms of the approximation and estimation error terms (up to constants and logarithmic factors), we obtain

$$\|\tilde{V} - V\|_\rho \leq O\left(\frac{1}{1-\gamma}\|V - \Pi V\|_\rho + \frac{1}{1-\gamma}\frac{1}{\sqrt{n}}\right).$$

While the first term (*approximation error*) only depends on the target function  $V$  and the function space  $\mathcal{F}$ , the second term (*estimation error*) primarily depends on the number of samples. Thus, when  $n$  goes to infinity, the estimation error goes to zero and we obtain the same performance bound (up to a  $4\sqrt{2}$  constant) as for the model-based case reported by Tsitsiklis & Van Roy (1997).

**Remark 2** Antos et al. (2008) reported a sample-based analysis for the modified Bellman residual (MBR) minimization algorithm. They consider a general setting in which the function space  $\mathcal{F}$  is bounded and the performance of the algorithm is evaluated according to an arbitrary measure  $\mu$  (possibly different than the stationary distribution of the Markov chain  $\rho$ ). Since Antos et al. (2008) showed that the MBR minimization algorithm is equivalent to LSTD when

$\mathcal{F}$  is a linearly parameterized space, it may be interesting to compare Theorem 2 to the bound in Lemma 11 of Antos et al. (2008). In Theorem 2, similar to Antos et al. (2008), samples are drawn from a stationary  $\beta$ -mixing process, however,  $\mathcal{F}$  is a linear space and  $\rho$  is the stationary distribution of the Markov chain. It is interesting to note the impact of these two differences in the final bound. The use of linear spaces has a direct effect on the estimation error and leads to a better convergence rate due to the use of improved functional concentration inequalities (Lemma 5 in the appendix). In fact, while in Antos et al. (2008) the estimation error for the squared error is of order  $O(1/\sqrt{n})$ , here we achieve a faster convergence rate of order  $O(1/n)$ . The use of  $\rho$  instead of an arbitrary measure  $\mu$  has a significant impact on the approximation error. The approximation error in Eq. 18  $\|V - \Pi V\|_\rho$  only depends on how well the space  $\mathcal{F}$  can approximate the value function  $V$ . On the other hand, the approximation error of Antos et al. (2008) contains terms that are related to more complex properties of the space, such as its capability to approximate any function obtained by applying the Bellman operator  $\mathcal{T}$  to any function in  $\mathcal{F}$ . This term called the inherent Bellman error can be shown to be small only for specific classes of MDPs (e.g., Lipschitz MDPs). Finally, it is interesting to notice that although the solution of MBR minimization reduces to LSTD, its sample-based analysis cannot be directly used for LSTD. In fact, in Antos et al. (2008) the function space  $\mathcal{F}$  is assumed to be bounded, while general linear spaces cannot be bounded. Whether the analysis of Antos et al. (2008) can be extended to the truncated solution  $\tilde{V}$  of LSTD is an open question that requires further investigation.

### 5.3. Generalization Bounds for Markov Chains

The main assumption in the previous section is that  $X_1, \dots, X_n$  is generated by a stationary  $\beta$ -mixing process with stationary distribution  $\rho$ . This is possible if we consider samples of a Markov chain during its stationary regime, i.e.  $X_1 \sim \rho$ . However in practice,  $\rho$  is not known, and the first sample  $X_1$  is usually drawn from a given initial distribution and the rest of the sequence is obtained by following the Markov chain from  $X_1$  on. As a result, the sequence  $X_1, \dots, X_n$  is no longer a realization of a stationary  $\beta$ -mixing process. Nonetheless, under suitable conditions, after  $\tilde{n} < n$  steps, the distribution of  $X_{\tilde{n}}$  approaches the stationary distribution  $\rho$ . In fact, according to the convergence theorem for fast-mixing Markov chains (see e.g., Proposition 3 in the appendix), for any initial distribution  $\lambda \in \mathcal{S}(\mathcal{X})$ , we have



$$\left\| \int_{\mathcal{X}} \lambda(dx) P^n(\cdot|x) - \rho(\cdot) \right\|_{TV} \leq \bar{\beta} \exp(-bn^\kappa).$$

We now derive a bound for a modification of pathwise LSTD in which the first  $\tilde{n}$  samples (that are used to burn the chain) are discarded and the remaining  $n - \tilde{n}$  samples are used as training samples for the algorithm.

**Theorem 3.** *Let  $X_1, \dots, X_n$  be a trajectory generated by a  $\beta$ -mixing Markov chain with stationary distribution  $\rho$ . Let  $\tilde{n}$  ( $1 \leq \tilde{n} < n$ ) be such that  $n - \tilde{n}$  satisfies the condition of Eq. 15, and  $X_{\tilde{n}+1}, \dots, X_n$  be the samples actually used by the algorithm. Let  $\omega$  be the smallest eigenvalue of the Gram matrix defined according to  $\rho$  and  $\alpha^* \in \mathbb{R}^d$  be such that  $f_{\alpha^*} = \Pi V$ . Let  $\tilde{V}$  be the truncation of the pathwise LSTD solution (using Eq. 14), then by setting  $\tilde{n} = \left(\frac{1}{b} \log \frac{2e\bar{\beta}n}{\delta}\right)^{1/\kappa}$ , we have*

$$\begin{aligned} \|\tilde{V} - V\|_\rho &\leq \frac{2}{1-\gamma} \left[ 2\sqrt{2}\|V - \Pi V\|_\rho + \varepsilon_2 \right. \\ &\quad \left. + \gamma V_{\max} L \sqrt{\frac{d}{\nu}} \left( \sqrt{\frac{8 \log(8d/\delta)}{n - \tilde{n}}} + \frac{1}{n - \tilde{n}} \right) \right] + \varepsilon_1 \end{aligned} \quad (19)$$

with probability  $1 - \delta$ , where  $\varepsilon_1$  and  $\varepsilon_2$  are defined as in Theorem 2 (with  $n - \tilde{n}$  as the number of training samples).

The proof of this result is a simple consequence of Lemma 8 in the appendix applied to Theorem 2.

**Remark 1** The bound in Eq. 19 indicates that in the case of  $\beta$ -mixing Markov chains, a similar performance to the one for stationary  $\beta$ -mixing processes is obtained by discarding the first  $\tilde{n} = O(\log n)$  samples.

## 6. Conclusions

In this paper we presented a finite-sample analysis of a natural version of LSTD, called pathwise LSTD. We first considered a general setting where we do not make any assumption about the Markov chain. We derived an empirical performance bound which indicates how close the LSTD solution is to the value function  $V$  at the states generated by the Markov chain. The bound is expressed in terms of the best possible approximation of  $V$  (approximation error) in the linear approximation space  $\mathcal{F}$  and an estimation error term which depends on the number of samples (the quadratic error scales with  $O(n^{-1/2})$ ) and the smallest strictly-positive eigenvalue of the sample-based Gram matrix. We then showed that when the Markov chain possesses a stationary distribution, then one can deduce generalization performance bounds using the stationary

distribution of the chain as our generalization measure. In particular, we considered the cases where the sample trajectory is generated by stationary and non-stationary  $\beta$ -mixing Markov chains and derived the corresponding bounds.

This work raises two open questions: 1) Is it possible to remove the dependence to  $\nu_n$  in the bound of Theorem 1? 2) Is it possible to extend the current analysis to the general case of LSTD( $\lambda$ )?

## 7. Appendix

### 7.1. IID Samples

Although in the setting considered in the paper the samples are non-i.i.d., we first report functional concentration inequalities for i.i.d. samples which will be later extended to stationary and non-stationary  $\beta$ -mixing processes.

We first recall the definition of empirical and expected norms for a function  $f$ :

$$\|f\|^2 = \mathbb{E}[|f(X_1)|^2], \quad \|f\|_{X_1^n}^2 = \frac{1}{n} \sum_{t=1}^n |f(X_t)|^2.$$

**Lemma 3.** *Let  $\mathcal{F}$  be a class of functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  bounded in absolute value by  $B$ . Let  $X_1^n = \{X_1, \dots, X_n\}$  be a sequence of i.i.d. samples. For any  $\epsilon > 0$*

$$\begin{aligned} \mathbb{P}[\exists f \in \mathcal{F} : \|f\| - 2\|f\|_{X_1^n} > \epsilon] \\ \leq 3\mathbb{E} \left[ \mathcal{N}_2 \left( \frac{\sqrt{2}}{24} \epsilon, \mathcal{F}, X_1^{2n} \right) \right] \exp \left( -\frac{n\epsilon^2}{288B^2} \right), \end{aligned} \quad (20)$$

and

$$\begin{aligned} \mathbb{P}[\exists f \in \mathcal{F} : \|f\|_{X_1^n} - 2\|f\| > \epsilon] \\ \leq 3\mathbb{E} \left[ \mathcal{N}_2 \left( \frac{\sqrt{2}}{24} \epsilon, \mathcal{F}, X_1^{2n} \right) \right] \exp \left( -\frac{n\epsilon^2}{288B^2} \right). \end{aligned} \quad (21)$$

*Proof.* The first statement was proved in Györfi et al. (2002) and the second one can be proved similarly.  $\square$

**Proposition 1.** *Let  $\mathcal{F}$  be a class of linear functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  of dimension  $d$  and  $\tilde{\mathcal{F}}$  be the class of functions obtained by truncating functions  $f \in \mathcal{F}$  at a threshold  $B$ . Then for any sample  $x_1^n$  and  $\epsilon > 0$*

$$\mathcal{N}_2 \left( \epsilon, \tilde{\mathcal{F}}, X_1^n \right) \leq 3 \left( \frac{3e(2B)^2}{\epsilon^2} \right)^{2(d+1)}. \quad (22)$$

*Proof.* Using Theorem 9.4. in Györfi et al. (2002) and the fact that the pseudo-dimension of  $\tilde{\mathcal{F}}$  is less than  $\mathcal{F}$ , we have

$$\begin{aligned} \mathcal{N}_2\left(\epsilon, \tilde{\mathcal{F}}, X_1^n\right) &\leq 3\left(\frac{2e(2B)^2}{\epsilon^2} \log \frac{3e(2B)^2}{\epsilon^2}\right)^{d+1} \\ &\leq 3\left(\frac{3e(2B)^2}{\epsilon^2}\right)^{2(d+1)}. \end{aligned}$$

□

We now use Proposition 1 to invert the bound in Lemma 3 for truncated linear spaces.

**Corollary 1.** *Let  $\mathcal{F}$  be a class of linear functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  of dimension  $d$ ,  $\tilde{\mathcal{F}}$  be the class of functions obtained by truncating functions  $f \in \mathcal{F}$  at a threshold  $B$ , and  $X_1^n = \{X_1, \dots, X_n\}$  be a sequence of i.i.d. samples. By inverting the bound of Lemma 3, for any  $\tilde{f} \in \tilde{\mathcal{F}}$ , we have*

$$\|\tilde{f}\| - 2\|\tilde{f}\|_{X_1^n} \leq \epsilon(\delta), \quad (23)$$

$$\|\tilde{f}\|_{X_1^n} - 2\|\tilde{f}\| \leq \epsilon(\delta), \quad (24)$$

with probability  $1 - \delta$ , where

$$\epsilon(\delta) = \sqrt{\frac{\Lambda(n, d, \delta)}{nC_2}}, \quad (25)$$

$\Lambda(n, d, \delta) = 2(d+1) \log n + \log \frac{\epsilon}{\delta} + \log(9(C_1 C_2)^{2(d+1)})$ ,  $C_1 = 3456eB^2$ , and  $C_2 = (288B^2)^{-1}$ .

*Proof.* In order to prove the corollary it is sufficient to verify that the following inequality holds for the  $\epsilon$  defined in (25)

$$3\mathbb{E}\left[\mathcal{N}_2\left(\frac{\sqrt{2}}{24}\epsilon, \tilde{\mathcal{F}}, X_1^{2n}\right)\right] \exp\left(-\frac{n\epsilon^2}{288B^2}\right) \leq \delta.$$

Using Proposition 1 we bound the first term as

$$\mathbb{E}\left[\mathcal{N}_2\left(\frac{\sqrt{2}}{24}\epsilon, \tilde{\mathcal{F}}, X_1^{2n}\right)\right] \leq 3\left(\frac{C_1}{\epsilon^2}\right)^{2(d+1)},$$

with  $C_1 = 3456eB^2$ . Next we notice that  $\Lambda(n, d, \delta) \geq 1$  and thus  $\epsilon \geq \sqrt{1/(nC_2)}$ . Using these bounds in the

original inequality and some algebra we obtain

$$\begin{aligned} &3\mathbb{E}\left[\mathcal{N}_2\left(\frac{\sqrt{2}}{24}\epsilon, \tilde{\mathcal{F}}, X_1^{2n}\right)\right] \exp\left(-\frac{n\epsilon^2}{288B^2}\right) \\ &\leq 9\left(\frac{C_1}{\epsilon^2}\right)^{2(d+1)} \exp(-nC_2\epsilon^2) \\ &\leq 9(nC_1 C_2)^{2(d+1)} \exp\left(-C_2 n \frac{\Lambda(n, d, \delta)}{nC_2}\right) \\ &= 9(nC_1 C_2)^{2(d+1)} n^{-2(d+1)} \frac{\delta}{e} \frac{1}{9(C_1 C_2)^{2(d+1)}} \\ &= \frac{\delta}{e} \leq \delta. \end{aligned}$$

□

Non-functional versions of Corollary 1 can be simply obtained by removing the covering number from the statement of Lemma 3.

**Corollary 2.** *Let  $f : \mathcal{X} \rightarrow \mathbb{R}$  be a bounded function and  $X_1^n = \{X_1, \dots, X_n\}$  be a sequence of i.i.d. samples. Then*

$$\|f\| - 2\|f\|_{X_1^n} \leq \epsilon(\delta), \quad (26)$$

$$\|f\|_{X_1^n} - 2\|f\| \leq \epsilon(\delta), \quad (27)$$

with probability  $1 - \delta$ , where

$$\epsilon(\delta) = \sqrt{\frac{288B^2}{n} \log \frac{3}{\delta}}. \quad (28)$$

## 7.2. Stationary $\beta$ -mixing Processes

We first introduce  $\beta$ -mixing stochastic processes and  $\beta$ -mixing coefficients.

**Definition 2.** *Let  $\{X_t\}_{t \geq 1}$  be a stochastic process. Let  $\sigma(X_i^j)$  denote the sigma-algebra generated by  $X_i^j$ . The  $m$ -th  $\beta$ -mixing coefficient of the stochastic process is defined by*

$$\beta_i = \sup_{t \geq 1} \mathbb{E} \left[ \sup_{B \in \sigma(X_{t+i}^\infty)} |P(B|X_1^t) - P(B)| \right].$$

*The process  $\{X_t\}_{t \geq 1}$  is said to be  $\beta$ -mixing if  $\beta_i \rightarrow 0$  as  $i \rightarrow \infty$ . In particular,  $\{X_t\}_{t \geq 1}$  mixes at an exponential rate with parameters  $\bar{\beta}, b, \kappa$  if  $\beta_i \leq \bar{\beta} \exp(-bi^\kappa)$ . Finally,  $\{X_t\}_{t \geq 1}$  is strictly stationary if  $X_t \sim \nu$  for any  $t > 0$ .*

Let  $X_1, \dots, X_n$  be a sequence of samples drawn from a stationary  $\beta$ -mixing process with coefficients  $\{\beta_i\}$ . We first introduce the blocking technique of Yu (1994).

Let us divide the sequence of samples into blocks of size  $k_n$ . For simplicity we assume  $n = 2m_n k_n$  with  $2m_n$  be the number of blocks.<sup>3</sup> For any  $1 \leq j \leq m_n$  we define the set of indexes in an odd and even block respectively as

$$H_j = \{t : 2(j-1)k_n + 1 \leq t \leq (2j-1)k_n\}$$

$$E_j = \{t : (2j-1)k_n + 1 \leq t \leq (2j)k_n\}$$

Let  $H = \cup_{j=1}^{m_n} H_j$  and  $E = \cup_{j=1}^{m_n} E_j$  be the set of all indexes in the odd and even blocks, respectively. We use  $X(H_j) = \{X_t : t \in H_j\}$  and  $X(H) = \{X_t : t \in H\}$ . We now introduce a ghost sample  $X'$  (the size of the ghost sample  $X'$  equals to the number of samples in each block  $k_n$ ) in each of the odd blocks such that the joint distribution of  $X'(H_j)$  is the same as  $X(H_j)$  but independent from any other block.

**Proposition 2.** (Yu, 1994) Let  $X_1, \dots, X_n$  be a sequence of samples drawn from a stationary  $\beta$ -mixing process with coefficients  $\{\beta_i\}$ . Let  $Q, Q'$  be the distributions of  $X(H)$  and  $X'(H)$ , respectively. For any measurable function  $h : \mathcal{X}^{m_n k_n} \rightarrow \mathbb{R}$  bounded by  $B$ , we have

$$|\mathbb{E}_Q[h(X(H))] - \mathbb{E}_{Q'}[h(X'(H))]| \leq B m_n \beta_{k_n}. \quad (29)$$

Before moving to the extension of Proposition 3 to  $\beta$  mixing processes, we report this technical lemma.

**Lemma 4.** Let  $a_1, \dots, a_n$  and  $b_1, \dots, b_n$  be two sequences of positive real numbers. Then

$$\left[ \left( \sum_{t=1}^n a_t^2 \right) \left( \sum_{t'=1}^n b_{t'}^2 \right) \right]^{1/2} \geq \sum_{t=1}^n a_t b_t \quad (30)$$

*Proof.* Using simple rules from algebra, we have

$$\begin{aligned} & \left[ \left( \sum_{t=1}^n a_t^2 \right) \left( \sum_{t'=1}^n b_{t'}^2 \right) \right]^{1/2} \\ &= \left[ \sum_{t=1}^n a_t^2 b_t^2 + \sum_{t=1}^n \sum_{\substack{t'=1 \\ t' \neq t}}^n a_t^2 b_{t'}^2 \right]^{1/2} \\ &= \left[ \left( \sum_{t=1}^n a_t b_t \right)^2 - 2 \sum_{t=1}^n \sum_{\substack{t'=1 \\ t' < t}}^n a_t b_t a_{t'} b_{t'} + \sum_{t=1}^n \sum_{\substack{t'=1 \\ t' \neq t}}^n a_{t'}^2 b_{t'}^2 \right]^{1/2} \\ &= \left[ \left( \sum_{t=1}^n a_t b_t \right)^2 + \sum_{t=1}^n \sum_{\substack{t'=1 \\ t' < t}}^n (a_t b_{t'} - a_{t'} b_t)^2 \right]^{1/2}. \end{aligned}$$

The claim follows because the second term in the bracket is always greater than or equal to zero.  $\square$

<sup>3</sup>The extension to the general case is straightforward.

**Lemma 5.** Let  $\mathcal{F}$  be a class of functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  bounded in absolute value by  $B$ ,  $X_1, \dots, X_n$  be a sequence of samples drawn from a stationary  $\beta$ -mixing process with coefficients  $\{\beta_i\}$ , and  $\epsilon > 0$ . Then we have

$$\mathbb{P}[\exists f \in \mathcal{F} : \|f\| - 2\|f\|_{X_1^n} > \epsilon] \leq 2\delta(\sqrt{2}\epsilon) + 2m_n \beta_{k_n}, \quad (31)$$

$$\mathbb{P}[\exists f \in \mathcal{F} : \|f\|_{X_1^n} - 2\sqrt{2}\|f\| > \epsilon] \leq 2\delta(\sqrt{2}\epsilon) + 2m_n \beta_{k_n}, \quad (32)$$

where

$$\delta(\epsilon) = 3\mathbb{E} \left[ \mathcal{N}_2 \left( \frac{\sqrt{2}}{24} \epsilon, \mathcal{F}, X(H) \cup X'(H) \right) \right] \exp \left( -\frac{m_n \epsilon^2}{288 B^2} \right).$$

*Proof.* Similar to Meir (2000), we first introduce  $\bar{\mathcal{F}}$  as the class of block functions  $\bar{f} : \mathcal{X}^{k_n} \rightarrow \mathbb{R}$  defined as

$$\bar{f}(X(H_j))^2 = \frac{1}{k_n} \sum_{t \in H_j} f(X_t)^2.$$

It is interesting to notice that block functions have exactly the same norms as the functions in  $\mathcal{F}$ . In fact

$$\begin{aligned} \|\bar{f}\|_{X(H)}^2 &= \frac{1}{m_n} \sum_{j=1}^{m_n} |\bar{f}(X(H_j))|^2 \\ &= \frac{1}{m_n} \sum_{j=1}^{m_n} \frac{1}{k_n} \sum_{t \in H_j} |f(X_t)|^2 = \|f\|_{X(H)}, \end{aligned} \quad (33)$$

and

$$\begin{aligned} \|\bar{f}\|^2 &= \mathbb{E} [|\bar{f}(X(H_1))|^2] \\ &= \frac{1}{k_n} \sum_{t \in H_1} \mathbb{E} [|f(X_t)|^2] = \mathbb{E} [|f(X_1)|^2] = \|f\|, \end{aligned} \quad (34)$$

where in Eq. 34, we used the fact that the process is stationary.

We now focus on Eq. 31

$$\begin{aligned} & \mathbb{P}[\exists f \in \mathcal{F} : \|f\| - 2\|f\|_{X_1^n} > \epsilon] \\ & \stackrel{(a)}{\leq} \mathbb{P}[\exists f \in \mathcal{F} : \|f\| - (\|f\|_{X(H)} + \|f\|_{X(E)}) > \epsilon] \\ & \stackrel{(b)}{=} \mathbb{P}[\exists f \in \mathcal{F} : \frac{1}{2}(\|f\| - 2\|f\|_{X(H)}) \\ & \quad + \frac{1}{2}(\|f\| - 2\|f\|_{X(E)}) > \epsilon] \\ & \stackrel{(c)}{\leq} \mathbb{P}[\exists f \in \mathcal{F} : \|f\| - 2\|f\|_{X(H)} > 2\epsilon] \\ & \quad + \mathbb{P}[\exists f \in \mathcal{F} : \|f\| - 2\|f\|_{X(E)} > 2\epsilon] \\ & \stackrel{(d)}{=} 2\mathbb{P}[\exists \bar{f} \in \bar{\mathcal{F}} : \|\bar{f}\| - 2\|\bar{f}\|_{X(H)} > 2\epsilon] \\ & \stackrel{(e)}{\leq} 2(\mathbb{P}[\exists \bar{f} \in \bar{\mathcal{F}} : \|\bar{f}\| - 2\|\bar{f}\|_{X'(H)} > 2\epsilon] + m_n \beta_{k_n}) \\ & \stackrel{(f)}{\leq} 2\delta'(2\epsilon) + 2m_n \beta_{k_n}. \end{aligned}$$

(a) We used the inequality  $\sqrt{a+b} \geq \frac{1}{\sqrt{2}}(\sqrt{a} + \sqrt{b})$  to split the norm  $\|f\|_{X_1^n} \geq \frac{1}{2}(\|f\|_{X(H)} + \|f\|_{X(E)})$ .

(b) Algebra.

(c) Split the probability.

(d) (1) Since the process is stationary the distribution over the even blocks is the same as the distribution over the odd blocks. (2) From Eqs. 33 and 34.

(e) Using Proposition 2 with  $h$  equals to the indicator function of the event inside the bracket, and the fact that the indicator function is bounded by  $B = 1$  and its expected value is equal to the probability of the event.

(f) Lemma 3 on space  $\bar{\mathcal{F}}$  where

$$\delta'(\epsilon) = 3\mathbb{E} \left[ \mathcal{N}_2 \left( \frac{\sqrt{2}}{24} \epsilon, \bar{\mathcal{F}}, \{X(H_j), X'(H_j)\}_{j=1}^{m_n} \right) \right] \exp \left( -\frac{m_n \epsilon^2}{288B^2} \right).$$

Now we relate the  $\ell_2$ -covering number of  $\bar{\mathcal{F}}$  to the covering number of  $\mathcal{F}$ . Using the definition of  $\bar{f}$  we have

$$\begin{aligned} \|\bar{f} - \bar{g}\|_{X(H)}^2 &= \frac{1}{m_n} \sum_{j=1}^{m_n} \left( \bar{f}(X(H_j)) - \bar{g}(X(H_j)) \right)^2 \\ &= \frac{1}{m_n k_n} \sum_{j=1}^{m_n} \left[ \left( \sum_{t \in H_j} f(X_t)^2 \right)^{\frac{1}{2}} - \left( \sum_{t' \in H_j} g(X_{t'})^2 \right)^{\frac{1}{2}} \right]^2. \end{aligned}$$

Taking the square and using Lemma 4, each element of the outer summation can be written as

$$\begin{aligned} &\sum_{t \in H_j} (f(X_t)^2 + g(X_t)^2) - 2 \left( \sum_{t \in H_j} f(X_t)^2 \right)^{\frac{1}{2}} \left( \sum_{t' \in H_j} g(X_{t'})^2 \right)^{\frac{1}{2}} \\ &\leq \sum_{t \in H_j} (f(X_t)^2 + g(X_t)^2 - 2f(X_t)g(X_t)) = \sum_{t \in H_j} (f(X_t) - g(X_t))^2. \end{aligned}$$

By taking the sum over all the odd blocks we obtain

$$\|\bar{f} - \bar{g}\|_{X(H)}^2 \leq \|f - g\|_{X(H)}^2,$$

which indicates that  $\mathcal{N}_2(\epsilon, \bar{\mathcal{F}}, \{X(H_j), X'(H_j)\}_{j=1}^{m_n}) \leq \mathcal{N}_2(\epsilon, \mathcal{F}, X(H) \cup X'(H))$ . Therefore, we have  $\delta'(2\epsilon) \leq \delta(2\epsilon) \leq \delta(\sqrt{2}\epsilon)$ , which concludes the proof.

With a similar approach, we can prove Eq. 32

$$\begin{aligned} &\mathbb{P} \left[ \exists f \in \mathcal{F} : \|f\|_{X_1^n} - 2\sqrt{2}\|f\| > \epsilon \right] \\ &\stackrel{(a)}{\leq} \mathbb{P} \left[ \exists f \in \mathcal{F} : \frac{\sqrt{2}}{2} (\|f\|_{X(H)} + \|f\|_{X(E)}) - 2\sqrt{2}\|f\| > \epsilon \right] \\ &\stackrel{(b)}{=} \mathbb{P} \left[ \exists f \in \mathcal{F} : \left( \frac{\sqrt{2}}{2} \|f\|_{X(H)} - \sqrt{2}\|f\| \right) \right. \\ &\quad \left. + \left( \frac{\sqrt{2}}{2} \|f\|_{X(E)} - \sqrt{2}\|f\| \right) > \epsilon \right] \\ &\stackrel{(c)}{\leq} \mathbb{P} \left[ \exists f \in \mathcal{F} : \|f\|_{X(H)} - 2\|f\| > \sqrt{2}\epsilon \right] \\ &\quad + \mathbb{P} \left[ \exists f \in \mathcal{F} : \|f\|_{X(E)} - 2\|f\| > \sqrt{2}\epsilon \right] \\ &\stackrel{(d)}{=} 2\mathbb{P} \left[ \exists \bar{f} \in \bar{\mathcal{F}} : \|\bar{f}\|_{X(H)} - 2\|\bar{f}\| > \sqrt{2}\epsilon \right] \\ &\stackrel{(e)}{\leq} 2 \left( \mathbb{P} \left[ \exists \bar{f} \in \bar{\mathcal{F}} : \|\bar{f}\|_{X'(H)} - 2\|\bar{f}\| > \sqrt{2}\epsilon \right] + m_n \beta_{k_n} \right) \\ &\stackrel{(f)}{\leq} 2\delta'(\sqrt{2}\epsilon) + 2m_n \beta_{k_n} \leq 2\delta(\sqrt{2}\epsilon) + 2m_n \beta_{k_n}. \end{aligned}$$

(a) We used the inequality  $\sqrt{a+b} \leq (\sqrt{a} + \sqrt{b})$  to split the norm  $\|f\|_{X_1^n} \leq \frac{\sqrt{2}}{2}(\|f\|_{X(H)} + \|f\|_{X(E)})$ .  $\square$

**Corollary 3.** Let  $\mathcal{F}$  be a class of linear functions  $f: \mathcal{X} \rightarrow \mathbb{R}$  of dimension  $d$ ,  $\bar{\mathcal{F}}$  be the class of functions obtained by truncating functions  $f \in \mathcal{F}$  at a threshold  $B$ , and  $X_1^n = \{X_1, \dots, X_n\}$  be a sequence of samples drawn from a stationary  $\beta$ -mixing process with coefficients  $\{\beta_i\}$ . By inverting the bound of Lemma 5, for any  $\tilde{f} \in \bar{\mathcal{F}}$  we have

$$\|\tilde{f}\| - 2\|\tilde{f}\|_{X_1^n} \leq \epsilon(\delta), \quad (35)$$

$$\|\tilde{f}\|_{X_1^n} - 2\sqrt{2}\|\tilde{f}\| \leq \epsilon(\delta), \quad (36)$$

with probability  $1 - \delta$ , where

$$\epsilon(\delta) = \sqrt{\frac{\Lambda(n, d, \delta)}{nC_2} \max \left\{ \frac{\Lambda(n, d, \delta)}{b}, 1 \right\}^{1/\kappa}}, \quad (37)$$

$$\begin{aligned} \Lambda(n, d, \delta) &= 2(d+1) \log n + \log \frac{\epsilon}{\delta} + \log^+ (\max\{18(C_1 C_2)^{2(d+1)}, \bar{\beta}\}), \quad C_1 = 1728eB^2, \\ &\text{and } C_2 = (288B^2)^{-1}. \end{aligned}$$

*Proof.* In order to prove the statement, we need to verify that  $\epsilon$  in Eq. 37 satisfies

$$\begin{aligned} \delta' &= 6\mathbb{E} \left[ \mathcal{N}_2 \left( \frac{1}{12} \epsilon, \bar{\mathcal{F}}, X(H) \cup X'(H) \right) \right] \exp \left( -\frac{m_n \epsilon^2}{144B^2} \right) \\ &\quad + 2m_n \beta_{k_n} \leq \delta. \end{aligned}$$

Using Proposition 1 the covering number can be bounded by

$$\mathbb{E} \left[ \mathcal{N}_2 \left( \frac{1}{12} \epsilon, \tilde{\mathcal{F}}, X(H) \cup X'(H) \right) \right] \leq 3 \left( \frac{1728eB^2}{\epsilon^2} \right)^{2(d+1)}.$$

By recalling the definition of the  $\beta$ -coefficients  $\{\beta_i\}$  and  $k_n \geq 1$  we have

$$2m_n \beta_{k_n} \leq \frac{n}{k_n} \bar{\beta} \exp(-bk_n^\kappa) \leq n\bar{\beta} \exp(-bk_n^\kappa).$$

From the last two inequalities and the definitions of  $C_1$  and  $C_2$ , and by setting  $D = 2(d+1)$  we obtain

$$\delta' \leq 18 \left( \frac{C_1}{\epsilon^2} \right)^D \exp \left( -\frac{nC_2\epsilon^2}{k_n} \right) + n\bar{\beta} \exp(-bk_n^\kappa).$$

By equalizing the arguments of the two exponential we obtain the definition of  $k_n$  as

$$k_n = \left\lceil \left( \frac{nC_2\epsilon^2}{b} \right)^{\frac{1}{\kappa+1}} \right\rceil.$$

Furthermore, we have

$$k_n \geq \max \left\{ \left( \frac{nC_2\epsilon^2}{b} \right)^{\frac{1}{\kappa+1}}, 1 \right\}, \quad \frac{1}{k_n} \geq \min \left\{ \left( \frac{b}{nC_2\epsilon^2} \right)^{\frac{1}{\kappa+1}}, 1 \right\}.$$

Using the above inequalities, we can write  $\delta'$  as

$$\begin{aligned} \delta' &\leq 18 \left( \frac{C_1}{\epsilon^2} \right)^D \exp \left( -\min \left\{ \frac{b}{nC_2\epsilon^2}, 1 \right\}^{\frac{1}{\kappa+1}} nC_2\epsilon^2 \right) \\ &\quad + n\bar{\beta} \exp \left( -b \max \left\{ \frac{nC_2\epsilon^2}{b}, 1 \right\}^{\frac{\kappa}{\kappa+1}} \right). \end{aligned}$$

The objective now is to make the arguments of the two exponential equal. For the second argument we have

$$\begin{aligned} &b \max \left\{ \frac{nC_2\epsilon^2}{b}, 1 \right\}^{\frac{\kappa}{\kappa+1}} \\ &= b \max \left\{ \frac{nC_2\epsilon^2}{b}, 1 \right\} \min \left\{ \frac{b}{nC_2\epsilon^2}, 1 \right\}^{\frac{1}{\kappa+1}} \\ &\geq nC_2\epsilon^2 \min \left\{ \frac{b}{nC_2\epsilon^2}, 1 \right\}^{\frac{1}{\kappa+1}}. \end{aligned}$$

Thus

$$\delta' \leq \left( 18 \left( \frac{C_1}{\epsilon^2} \right)^D + n\bar{\beta} \right) \exp \left( -\min \left\{ \frac{b}{nC_2\epsilon^2}, 1 \right\}^{\frac{1}{\kappa+1}} nC_2\epsilon^2 \right).$$

Now we plug in  $\epsilon$  from Eq. 37. Using the fact that  $\Lambda \geq 1$ , we know that  $\epsilon^2 \geq (nC_2)^{-1}$ , and thus

$$\delta' \leq \left( 18 (nC_1C_2)^D + n\bar{\beta} \right) \exp(-\Lambda).$$

Using the definition of  $\Lambda$ , we obtain

$$\begin{aligned} \delta' &\leq \left( 18 (nC_1C_2)^D + n\bar{\beta} \right) n^{-D} \max\{18(C_1C_2)^D, \bar{\beta}\}^{-1} \frac{\delta}{e} \\ &\leq (1 + n^{1-D}) \frac{\delta}{e} \leq (1 + 1) \frac{\delta}{e} \leq \delta, \end{aligned}$$

which concludes the proof.  $\square$

In order to understand better the shape of the estimation error, we consider a simple  $\beta$ -mixing process with parameters  $\bar{\beta} = b = \kappa = 1$ . Eq. 37 reduces to

$$\epsilon(\delta) = \sqrt{\frac{288B^2\Lambda(n, d, \delta)^2}{n}},$$

with  $\Lambda(n, d, \delta) = 2(d+1) \log n + \log \frac{\epsilon}{\delta} + \log(18(6e)^{2(d+1)})$ .

Finally, we report the non-functional version of the previous corollary.

**Corollary 4.** *Let  $f : \mathcal{X} \rightarrow \mathbb{R}$  be a linear function,  $\tilde{f}$  be its truncation at a threshold  $B$ , and  $X_1^n = \{X_1, \dots, X_n\}$  be a sequence of samples drawn from a stationary  $\beta$ -mixing process with coefficients  $\{\beta_i\}$ . Then*

$$\|\tilde{f}\| - 2\|\tilde{f}\|_{X_1^n} \leq \epsilon(\delta), \quad (38)$$

$$\|\tilde{f}\|_{X_1^n} - 2\sqrt{2}\|\tilde{f}\| \leq \epsilon(\delta), \quad (39)$$

with probability  $1 - \delta$ , where

$$\epsilon(\delta) = \sqrt{\frac{\Lambda(n, \delta)}{nC_2} \max \left\{ \frac{\Lambda(n, \delta)}{b}, 1 \right\}^{1/\kappa}}, \quad (40)$$

$\Lambda(n, \delta) = \log \frac{\epsilon}{\delta} + \log(\max\{6, n\bar{\beta}\})$ , and  $C_2 = (288B^2)^{-1}$ .

*Proof.* The proof follows the same steps as in Corollary 3. We have the following sequence of inequalities

$$\begin{aligned} \delta' &\leq 6 \exp \left( -\frac{nC_2\epsilon^2}{k_n} \right) + \frac{n}{k_n} \bar{\beta} \exp(-bk_n^\kappa) \\ &\leq (6 + n\bar{\beta}) \exp(-\Lambda) \\ &= (6 + n\bar{\beta}) \max\{6, n\bar{\beta}\}^{-1} \frac{\delta}{e} \leq (1 + 1) \frac{\delta}{e} \leq \delta. \end{aligned}$$

$\square$

### 7.3. Markov Chains

We first review the conditions for the convergence of Markov chains (Theorem 13.3.3. in [Meyn & Tweedie 1993](#)).

**Proposition 3.** *Let  $\mathcal{M}$  be an ergodic and aperiodic Markov chain defined on  $\mathcal{X}$  with stationary distribution  $\rho$ . If  $P(A|x)$  is the transition kernel of  $\mathcal{M}$  with  $A \subseteq \mathcal{X}$  and  $x \in \mathcal{X}$ , then for any initial distribution  $\lambda$*

$$\lim_{i \rightarrow \infty} \left\| \int_{\mathcal{X}} \lambda(dx) P^i(\cdot|x) - \rho(\cdot) \right\|_{TV} = 0, \quad (41)$$

where  $\|\cdot\|_{TV}$  is the total variation norm.



**Definition 3.** Let  $\mathcal{M}$  be an ergodic and aperiodic Markov chain with stationary distribution  $\rho$ .  $\mathcal{M}$  is mixing with an exponential rate with parameters  $\bar{\beta}, b, \kappa$ , if its  $\beta$ -mixing coefficients  $\{\beta_i\}$  are defined as  $\beta_i \leq \bar{\beta} \exp(-bi^\kappa)$ . Then for any initial distribution  $\lambda$

$$\| \int_{\mathcal{X}} \lambda(dx) P^i(\cdot|x) - \rho(\cdot) \|_{TV} \leq \bar{\beta} \exp(-bi^\kappa). \quad (42)$$

**Lemma 6.** Let  $\mathcal{M}$  be an ergodic and aperiodic Markov chain with a stationary distribution  $\rho$ . Let  $X_1, \dots, X_n$  be a sequence of samples drawn from the stationary distribution of the Markov chain  $\rho$  and  $X'_1, \dots, X'_n$  be a sequence of samples such that  $X'_1 \sim \rho'$  and  $X'_{1 < t \leq n}$  are generated by simulating  $\mathcal{M}$  from  $X'_1$ . Let  $\eta$  be an event defined on  $\mathcal{X}^n$ , then

$$|\mathbb{P}[\eta(X_1, \dots, X_n)] - \mathbb{P}[\eta(X'_1, \dots, X'_n)]| \leq \|\rho' - \rho\|_{TV} \quad (43)$$

*Proof.* We prove one side of the inequality. Let  $Q$  be the conditional joint distribution of  $(X_{1 < t \leq n} | X_1 = x)$  and  $Q'$  be the conditional joint distribution of  $(X'_{1 < t \leq n} | X'_1 = x)$ . We first notice that  $Q$  is exactly the same as  $Q'$ . In fact, the first sequence  $(X_{1 < t \leq n})$  is generated by drawing  $X_1$  from the stationary distribution  $\rho$  and then following the Markov chain. Similarly, the second sequence  $(X'_{1 < t \leq n})$  is obtained following the Markov chain from  $X'_1 \sim \rho'$ . As a result, the conditional distributions of the two sequences is exactly the same and just depend on the Markov chain. As a result, we obtain the following sequence of inequalities

$$\begin{aligned} & \mathbb{P}[\eta(X_1, \dots, X_n)] \\ &= \mathbb{E}_{X_1, \dots, X_n} [\mathbb{I}\{\eta(X_1, \dots, X_n)\}] \\ &= \mathbb{E}_{X_1 \sim \rho} [\mathbb{E}_{X_2, \dots, X_n} [\mathbb{I}\{\eta(X_1, X_2, \dots, X_n)\} | X_1]] \\ &= \mathbb{E}_{X_1 \sim \rho} [\mathbb{E}_{X'_2, \dots, X'_n} [\mathbb{I}\{\eta(X_1, X'_2, \dots, X'_n)\} | X_1]] \\ &\stackrel{(a)}{\leq} \mathbb{E}_{X_1 \sim \rho'} [\mathbb{E}_{X'_2, \dots, X'_n} [\mathbb{I}\{\eta(X_1, X'_2, \dots, X'_n)\} | X_1]] \\ &\quad + \|\rho' - \rho\|_{TV} \\ &\stackrel{(b)}{=} \mathbb{E}_{X'_1 \sim \rho'} [\mathbb{E}_{X'_2, \dots, X'_n} [\mathbb{I}\{\eta(X'_1, X'_2, \dots, X'_n)\} | X'_1]] \\ &\quad + \|\rho' - \rho\|_{TV} \\ &= \mathbb{P}[\eta(X'_1, \dots, X'_n)] + \|\rho' - \rho\|_{TV}. \end{aligned}$$

Note that  $\mathbb{I}\{\cdot\}$  is the indicator function.

(a) simply follows from

$$\begin{aligned} & \mathbb{E}_{X \sim \rho} [f(X)] - \mathbb{E}_{X \sim \rho'} [f(X)] \\ &= \int_{\mathcal{X}} f(x) \rho(dx) - \int_{\mathcal{X}} f(x) \rho'(dx) \\ &\leq \|f\|_{\infty} \int_{\mathcal{X}} (\rho(dx) - \rho'(dx)) \leq \|f\|_{\infty} \|\rho - \rho'\|_{TV}. \end{aligned}$$

(b) From the fact that  $X_1 = X'_1 = x$ . □

**Lemma 7.** Let  $\mathcal{F}$  be a class of functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  bounded in absolute value by  $B$ ,  $\mathcal{M}$  be an ergodic and aperiodic Markov chain with a stationary distribution  $\rho$ . Let  $\mathcal{M}$  be mixing with an exponential rate with parameters  $\bar{\beta}, b, \kappa$ . Let  $\lambda$  be an initial distribution over  $\mathcal{X}$  and  $X_1, \dots, X_n$  be a sequence of samples such that  $X_1 \sim \lambda$  and  $X_{1 < t \leq n}$  obtained by following  $\mathcal{M}$  from  $X_1$ . For any  $\epsilon > 0$ ,

$$\begin{aligned} & \mathbb{P}[\exists f \in \mathcal{F} : \|f\| - 2\|f\|_{X_1^n} > \epsilon] \\ & \leq \|\lambda - \rho\|_{TV} + 2\delta(\sqrt{2}\epsilon) + 2m_n\beta_{k_n}, \end{aligned}$$

and

$$\begin{aligned} & \mathbb{P}[\exists f \in \mathcal{F} : \|f\|_{X_1^n} - 2\sqrt{2}\|f\| > \epsilon] \\ & \leq \|\lambda - \rho\|_{TV} + 2\delta(\sqrt{2}\epsilon) + 2m_n\beta_{k_n}, \end{aligned}$$

where

$$\delta(\epsilon) = 3\mathbb{E} \left[ \mathcal{N}_2 \left( \frac{\sqrt{2}}{24}\epsilon, \mathcal{F}, X(H) \cup X'(H) \right) \right] \exp \left( -\frac{m_n\epsilon^2}{288B^2} \right).$$

*Proof.* The proof is an immediate consequence of Lemma 5 and Lemma 6 by defining  $\eta(X_1, \dots, X_n)$  as

$$\eta(X_1, \dots, X_n) = \{\exists f \in \mathcal{F} : \|f\| - 2\|f\|_{X_1^n} > \epsilon\},$$

and

$$\eta(X_1, \dots, X_n) = \{\exists f \in \mathcal{F} : \|f\|_{X_1^n} - 2\sqrt{2}\|f\| > \epsilon\},$$

respectively. □

Finally, we consider a special case in which out of the  $n$  total number of samples,  $\tilde{n}$  ( $1 \leq \tilde{n} < n$ ) are used to “burn” the chain and  $n - \tilde{n}$  are actually used as training samples.

**Lemma 8.** Let  $\mathcal{F}$  be a class of linear functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  of dimension  $d$  and  $\tilde{\mathcal{F}}$  be the class of functions obtained by truncating functions  $f \in \mathcal{F}$  at a threshold  $B$ . Let  $\mathcal{M}$  be an ergodic and aperiodic Markov chain with a stationary distribution  $\rho$ . Let  $\mathcal{M}$  be mixing with an exponential rate with parameters  $\bar{\beta}, b, \kappa$ . Let  $\mu$  be the initial distribution and  $X_1, \dots, X_n$  be a sequence of samples such that  $X_1 \sim \mu$  and  $X_{1 < t \leq n}$  obtained by following  $\mathcal{M}$  from  $X_1$ . If the first  $\tilde{n}$  ( $1 \leq \tilde{n} < n$ ) samples are used to burn the chain and  $n - \tilde{n}$  are actually used as training samples, by inverting Lemma 7, for any  $\tilde{f} \in \tilde{\mathcal{F}}$ , we obtain

$$\|\tilde{f}\| - 2\|\tilde{f}\|_{X_1^n} \leq \epsilon(\delta),$$

$$\|\tilde{f}\|_{X_1^n} - 2\sqrt{2}\|\tilde{f}\| \leq \epsilon(\delta),$$

with probability  $1 - \delta$ , where

$$\epsilon(\delta) = \sqrt{\frac{\Lambda(n - \tilde{n}, d, \delta)}{(n - \tilde{n})C_2} \max \left\{ \frac{\Lambda(n - \tilde{n}, d, \delta)}{b}, 1 \right\}^{1/\kappa}},$$

$$\Lambda(n, d, \delta) = 2(d + 1) \log n + \log \frac{\epsilon}{\delta} + \log^+ (\max\{18(C_1 C_2)^{2(d+1)}, \bar{\beta}\}), \quad C_1 = 1728eB^2,$$

$$C_2 = (288B^2)^{-1}, \text{ and } \tilde{n} = \left( \frac{1}{b} \log \frac{2e\bar{\beta}n}{\delta} \right)^{1/\kappa}.$$

*Proof.* After  $\tilde{n}$  steps, the first sample used in the training set  $(X_{\tilde{n}+1})$  is drawn from the distribution  $\lambda = \mu P^{\tilde{n}}$ . Using Proposition 3 and Definition 3 we have

$$\|\lambda - \rho\|_{TV} \leq \bar{\beta} \exp(-b\tilde{n}^\kappa). \quad (44)$$

We first substitute the total variation in Lemma 7 with the bound in Eq. 44, and then verify that  $\epsilon$  in Eq. 8 satisfies the following inequality.

$$\begin{aligned} \delta' &= \|\lambda - \rho\|_{TV} + 2\delta(\sqrt{2}\epsilon) + 2m_{n-\tilde{n}}\beta_{k_{n-\tilde{n}}} \\ &\leq \bar{\beta} \exp(-b\tilde{n}^\kappa) \\ &\quad + 18 \left( \frac{C_1}{\epsilon^2} \right)^D \exp \left( -\frac{(n - \tilde{n})C_2\epsilon^2}{k_{n-\tilde{n}}} \right) \\ &\quad + (n - \tilde{n})\bar{\beta} \exp(-bk_{n-\tilde{n}}^\kappa) \\ &\leq \left( \frac{1}{2n} + 1 + (n - \tilde{n})^{1-D} \right) \frac{\delta}{e} \leq \left( \frac{1}{2} + 1 + 1 \right) \frac{\delta}{e} \\ &\leq \delta \end{aligned}$$

The above inequality can be verified by following the same steps as in Corollary 3 and by optimizing the bound for  $\tilde{n}$ .  $\square$

## References

- Antos, A., Szepesvári, Cs., and Munos, R. Learning near-optimal policies with Bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning Journal*, 71:89–129, 2008.
- Bertsekas, D. *Dynamic Programming and Optimal Control*. Athena Scientific, 2001.
- Boyan, J. Least-squares temporal difference learning. *Proceedings of the 16th International Conference on Machine Learning*, pp. 49–56, 1999.
- Bradtke, S. and Barto, A. Linear least-squares algorithms for temporal difference learning. *Machine Learning*, 22:33–57, 1996.
- Györfi, L., Kohler, M., Krzyżak, A., and Walk, H. *A distribution-free theory of nonparametric regression*. Springer-Verlag, New York, 2002.
- Lagoudakis, M. and Parr, R. Least-squares policy iteration. *Journal of Machine Learning Research*, 4: 1107–1149, 2003.
- Meir, R. Nonparametric time series prediction through adaptive model selection. *Machine Learning*, 39(1): 5–34, April 2000.
- Meyn, S. P. and Tweedie, R. L. *Markov chains and stochastic stability*. Springer-Verlag, 1993.
- Sutton, R. and Barto, A. *Reinforcement Learning: An Introduction*. MIP Press, 1998.
- Tsitsiklis, J. and Van Roy, B. An analysis of temporal difference learning with function approximation. *IEEE Transactions on Automatic Control*, 42:674–690, 1997.
- Yu, B. Rates of convergence for empirical processes of stationary mixing sequences. *The Annals of Probability*, 22(1):94–116, January 1994.